

# Synonymous mutations reduce genome compactness in icosahedral ssRNA viruses

Luca Tubiana,<sup>1,\*</sup> Anže Lošdorfer Božič,<sup>1,2</sup> Cristian Micheletti,<sup>3</sup> and Rudolf Podgornik<sup>1,4,5</sup>

<sup>1</sup>*Department of Theoretical Physics, Jožef Stefan Institute, SI-1000 Ljubljana, Slovenia*

<sup>2</sup>*Max Planck Institute for Biology of Ageing, Joseph-Stelzmann-Str. 9b, D-50931 Cologne, Germany*

<sup>3</sup>*SISSA, Via Bonomea 265, I-34136 Trieste, Italy*

<sup>4</sup>*Department of Physics, Faculty of Mathematics and Physics,  
University of Ljubljana, SI-1000 Ljubljana, Slovenia*

<sup>5</sup>*Department of Physics, University of Massachusetts, Amherst, MA 01003*

(Dated: October 30, 2014)

Recent studies have shown that single-stranded viral RNAs fold into more compact structures than random RNA sequences with similar chemical composition and identical length. Based on this comparison it has been suggested that wild-type viral RNA may have evolved to be atypically compact so as to aid its encapsidation and assist the viral assembly process. In order to further explore the compactness selection hypothesis, we systematically compare the predicted sizes of more than one hundred wild-type viral sequences with those of their mutants, which are evolved *in silico* and subject to a number of known evolutionary constraints. In particular, we enforce mutation synonymity, preserve the codon-bias, and leave untranslated regions intact. It is found that progressive accumulation of these restricted mutations still suffices to completely erase the characteristic compactness imprint of the viral RNA genomes, making them in this respect physically indistinguishable from randomly shuffled RNAs. This shows that maintaining the physical compactness of the genome is indeed a primary factor among ssRNA viruses evolutionary constraints, contributing also to the evidence that synonymous mutations in viral ssRNA genomes are not strictly neutral.

---

\* luca.tubiana@ijs.si; Corresponding author

## INTRODUCTION

Minimalistic organisms, such as single-stranded (ss)RNA viruses, are ideally suited to investigate how the three-dimensional organization of the genome – and not only its sequence composition – is subject to selective evolutionary pressure. We recall, for instance, that several structural features are robustly maintained in the highly-mutating ssRNA viruses. These include RNA structures acting as signals for translation [1], for transcription initiation [2], or as packaging signals to initiate the self-assembly of the virion [3, 4]. Other conserved structures have also been identified [5–7], including long-range interactions between different genomic regions of RNA [5, 8], whose role in the virus life cycle is still unknown.

The preservation of these structural features must act as a powerful constraint on viable RNAs, together with the multiple other, often competing, selection pressures [9–11]. The evolutionary mechanisms which maintain the viral protein phenotype clearly impact the genome chemical composition more directly, by largely restricting those mutations which have a deleterious effect on the encoded proteins [12–15]. On the other hand, synonymous mutations, i.e., mutations that do not change the amino acid sequences encoded by the genes, are neutral with regard to these mechanisms, but still have an impact on the structural features of RNAs.

It is increasingly becoming recognised that the mechanisms which may constrain synonymous mutations extend beyond the aforementioned conservation of specific genome structures, and are underpinned by general physico-chemical constraints. The latter mostly stem from the polymeric nature of the gene-carrying macromolecules and their steric and electrostatic self-interactions, as well as interactions with the capsid proteins [16–19]. These molecular interactions can be long-ranged and depend crucially on the pH of the local aqueous solution environment [20], conferring virions the ability to assemble and disassemble spontaneously at proper bathing solution conditions [21–28], and the ability to recognize and selectively encapsidate only viral RNA even in the absence of packaging signals [19, 29–32].

In this study we focus on a general and major structure-related selection constraint, namely the feasibility to efficiently package viral RNA inside the capsid, and address its competition with sequence-based selection mechanisms. The overarching question is whether the viral RNA sequence has evolved not only for encoding a specific protein phenotype but also for promoting an innate fold of the free (unencapsidated) viral RNA itself that is primed for efficient encapsidation.

Major advances towards solving this important conundrum have been recently made by comparing the predicted equilibrium properties of ssRNA folds of several icosahedral viruses with those of random RNA sequences with similar length and nucleotide composition. By using general arguments based on the scaling properties of linear [33] and/or branched polymers [34], the folded wild-type (WT) viral RNA was shown to be significantly more compact than random nucleotide sequences. In addition – and most notably – the average radius of gyration of WT RNA genomes was found to exceed only slightly the inner radius of the fully-assembled capsid [35].

In this context, a key and still open problem relates to the extent to which the selective pressure for easily encapsidable RNA genomes directly competes with the other sequence-based mechanisms that are simultaneously at play for selecting biologically-viable viral RNA. As a matter of fact, the enhanced compactness of viral RNA has so far been established only by comparison against random sequences that do not retain any specific viral-like characteristics except from the overall nucleotide composition. As the volume of the sequence “phase space” that is accessible to viable viral RNA sequences is actually vanishingly small compared to the available combinatorial phase space of random sequences, it is crucial to ascertain the implications of introducing realistic sequence constraints into the picture. Such constraints could even affect the properties of the associated folds to the point of implying genome compactness, which would make the assumption of a distinct selection principle based on RNA compactness superfluous.

To address these issues we consider the implications of constrained mutations that conserve the encoded protein phenotype and the viral-like nucleotide composition on the compactness of viral RNA genomes. This allows us to examine the concurrence, or possibly the incompatibility, of sequence- and structure-based *parallel selection mechanisms*, and to ascertain whether the conservation of RNA compactness is among the causes of the sensitivity of ssRNA viruses to synonymous mutations.

Specifically, we consider 128 viral RNA sequences and evolve them synthetically by accumulating exclusively *synonymous* point-wise mutations, measuring their impact on the properly quantified compactness of the genome. We recall that the constraint of synonymity, i.e., considering only codons that encode for the same amino acids, is particularly severe for viral RNA because of both the high gene density and the frequent presence of overlapping reading frames.

Our study unequivocally shows that, at least for the viruses studied, the accumulation of strictly synonymous mutations – even if they are sparse – is sufficient to cause a systematic drift of the properly quantified compactness of the genome towards values comparable to those of unrestricted random sequences that are systematically much larger than those of the WT genomes. By focusing on the mutational dynamics of four viral genomes we show that while mutating as few as 5 % of a genome is enough to erase its compactness, there is still a non-negligible portion of the

sequence space in the vicinity of the WT sequence in which the genomes are at least as compact as the WT genome, while still coding for the correct proteins.

Furthermore, we show that the typical WT RNA compactness is related neither to the codon usage biases present in viral genomes nor to the particular sequences of the untranslated regions (UTRs) present at the 5' and 3' ends of the genomes. These results provide *a posteriori* evidence that the same viral RNA sequence can encode not only for the expression of the proper protein complement, exposed to canonical selection pressure mechanisms, but can on another level also prime the optimal physico-chemical genome-packing organization.

## MATERIALS AND METHODS

### Wild-type viral sequences

Viral ssRNA sequences were obtained from the NCBI nucleotide database [36]. The dataset we use includes positive-strand ssRNA viruses from the following families: Tymoviridae (from the order Tymovirales), Flaviviridae, Caliciviridae, Picornaviridae, Comovirinae, Bromoviridae, and Tombusviridae [37]. All the viruses considered have icosahedral capsids, the majority of them with triangulation number  $T = 3$ . Most of the families in the dataset have monopartite genome, with the exception of Comovirinae, which have a bipartite genome, and Bromoviridae, which have a tripartite genome [37]. Comovirinae pack the two segments, denoted RNA1 and RNA2, into separate virions; the two largest RNA segments of Bromoviridae genome, denoted RNA1 and RNA2, are also packed into separate virions, and we thus consider only these two segments. All the considered viruses use the eukaryotic genetic code and their genes have no reading gaps. Several sequences among those we consider also have overlapping reading frames, which are known to impose further evolutionary constraints increasing the deleterious effects of mutations [38, 39]. With these restrictions taken into account, the final dataset of analyzed sequences contains 128 viral genomes (compiled in Supporting Information (SI) Table S1).

### Synonymous point mutations

Extended models of sequence evolution of overlapping genes can account for the codependency of the nucleotide substitution process in two reading frames [40, 41], but are based in computationally very intense simulations and are not always applicable to large sequence datasets. Since in the present study we are interested in the statistical properties across various viral families, we adopt a much simpler model which simply conserves the produced amino acids in all reading frames.

Mutated viral ssRNA sequences are obtained using a Monte Carlo (MC) scheme designed to simulate synonymous point substitutions while also conserving dinucleotide frequencies. Starting from a WT sequence, a point substitution is introduced at every step and accepted or rejected using a Metropolis algorithm. Substitutions which change the amino acids encoded by the genes and are thus non-synonymous are rejected. To preserve the dinucleotide frequencies we additionally introduce a fictitious energy related to the viral dinucleotide odd-ratios [42]:

$$E = \sum_{XY} K_{XY} [O(XY) - O_{WT}(XY)]^2, \quad (1)$$

where

$$O(XY) = \frac{N(XY)}{N(X)N(Y)}N, \quad X, Y \in \{A, U, G, C\}. \quad (2)$$

Here,  $N(XY)$  is the number of  $XY$  pairs,  $N(X), N(Y)$  are the numbers of  $X$  and  $Y$  nucleotides in the sequence, and  $N$  is the total length of the RNA sequence.

The values of the constants  $K_{XY}$  are chosen in such a way that a considerable portion (but not all) of the proposed sequences have dinucleotide odd-ratios lying within  $1.5\Delta Q$ , where  $\Delta Q$  is the interquartile distance evaluated on the  $O_{WT}(XY)$  distribution of the corresponding viral family (see SI for additional information). We produce an extensive ensemble of point mutations ( $\sim 10^9$ ) to ensure an appropriate sampling of the sequence space. Sequences are sampled every  $100N$  mutations to ensure they are uncorrelated, and filtered *a posteriori* to have all odd-ratios within  $1.5\Delta Q$ . For every WT viral sequence we generate a set of 500 to 2000 mutated sequences and finally characterise the spatial compactness of the associated fold by computing the thermally averaged maximum ladder distance,  $\langle \text{MLD} \rangle$ , described in a later subsection.

As an additional check we also produce synonymous substitutions using the Fisher-Yates shuffling algorithm [43, 44] – in this way, the exact chemical composition of the sequences is conserved, although the dinucleotide odd-ratios are not. While much more complex models for the nucleotide substitutions exist (see for instance the review by Anisimova and Kosiol [45] and references therein), we chose these simple ones that conserve the chemical composition of the sequences as they are sufficient to prove our point, and can most importantly be applied in the same manner to all the genomes we considered.

To investigate the effect of progressively accumulating mutations on viral RNA compactness, quantified by the MLD, we first choose the  $K_{XY}$  values in such a way that all produced sequences obey the dinucleotide constraints. The generated MC trajectories are then sampled every  $N/100$  steps. This sampling produces strongly correlated sequences which show the evolution of the genome MLDs toward the values of their random counterparts.

### Synonymous mutations preserving codon bias

As an optional additional constraint, we fix the WT codon population by shuffling equivalent codons, as done in Ref. [46]. The shuffling is performed at the gene-wise level by first enumerating and pooling the synonymous codons in the WT gene sequence. Each codon in the latter is then replaced by one picked randomly from its synonymous pool. The pools are thus progressively depleted until all reassignments are completed, as in the standard Fisher-Yates shuffling algorithm [43, 44]. This shuffling procedure, which clearly preserves the WT codon bias at the gene level is applicable to viral genomes without overlapping genes, which are 86 in our case.

### Random RNA sequences

Random ssRNA sequences, used to obtain the scaling law for the MLD of random RNAs, are produced by shuffling RNA sequences with the Fisher-Yates algorithm [43, 44]. Random numbers, here as well as in the rest of the paper, are generated by the SIMD-oriented Fast Mersenne Twister (SFMT) random generator, version 1.4 [47]. The SFMT has a period of  $2^{216091} - 1$ , which suffices to produce random permutations of even 10 knt long RNA sequences. We use the same viral-like composition for the random sequences as in Ref. [33], that is, 0.26 A, 0.28 U, 0.24 G, and 0.22 C, to obtain the scaling law for random viral-like RNAs. This average composition is computed excluding Tymoviridae, which differ significantly in their composition. For the Tymoviridae family, we use the averaged composition of the viruses in our sample belonging to this family only (see SI Table S1 for the list), with the corresponding nucleotide composition: 0.219 A, 0.254 U, 0.163 G, 0.364 C.

### Maximum Ladder Distance (MLD)

In order to investigate the possibility that synonymous substitutions, while being neutral with respect to the encoded protein complement, can affect the secondary structure of viral RNA, we use the (thermally averaged) MLD, a quantitative, albeit coarse-grained indicator of the compactness of RNA folds introduced by Yoffe and coworkers [33]. While the MLD of random RNAs with viral-like nucleotide composition follows a simple scaling law, the MLDs of viral ssRNA genomes are on the other hand significantly lower, indicating that their folds are more compact than those of random RNAs.

When treating the RNA as an ideal linear polymer, one can compute its MLD when mapped to an ideal graph [33, 48]: For every pair of nucleotides  $i$  and  $j$  in an RNA sequence we compute the ladder distance, i.e., the number of steps on the ladder which separates the two nucleotides on the folded RNA. The maximum value of all the ladder distances in a fold is then its MLD; an example is shown in Fig. 1(a). By treating the MLD contour as the backbone of a linear polymer chain, this provides a measure of compactness/extendedness of the RNA molecule, even though it is not a direct measure of the three-dimensional size of the RNA. This simple measure yields the same scaling relationships as in the case when one treats the RNA as an ideal branched polymer, computing its root-mean-square radius of gyration to determine its extendedness [34].

The secondary structures of viral and random RNA sequences for which we determine their MLDs are obtained by folding the sequences with the `RNAsubopt` program available in the ViennaRNA Package, version 2.1 [49]. Due to the length of viral RNA, a population of different folds having comparable energy is expected. Therefore, instead of looking for the minimum energy fold, we produce 500 folds at thermal equilibrium for every RNA sequence. This results in a thermal average for the MLD of every sequence, obtained by averaging over this ensemble.

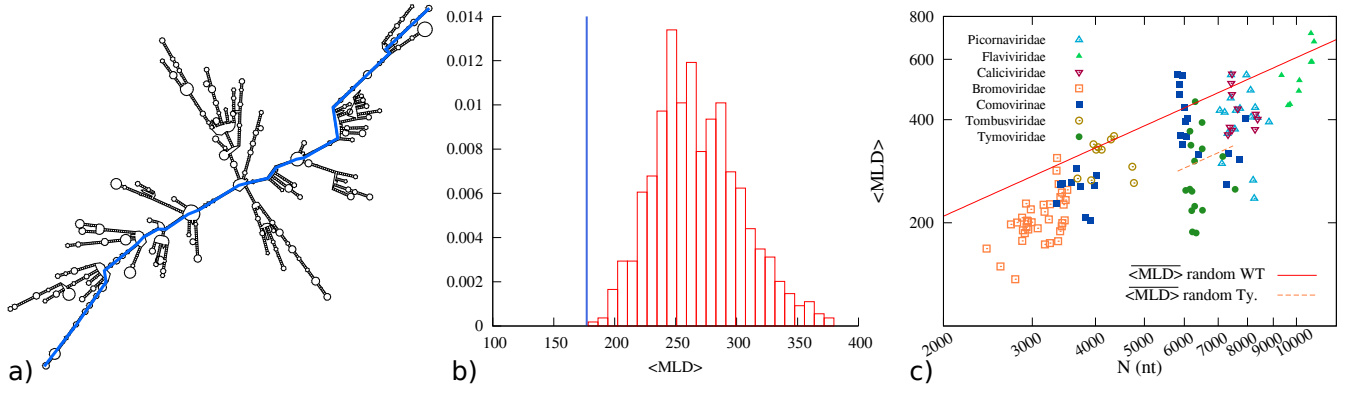


Figure 1. a) Example of a typical fold of the entire brome mosaic virus (BMV) RNA2 sequence. The maximum ladder distance (MLD) of the folded sequence is highlighted. b) Thermally averaged MLD,  $\langle \text{MLD} \rangle$ , of the WT BMV RNA2 sequence (blue line) and the distribution of  $\langle \text{MLD} \rangle$ s obtained for random RNA sequences of same length and composition as the WT sequence. c)  $\langle \text{MLD} \rangle$  of viral ssRNA sequences versus the sequence length  $N$  (in nucleotides). Different virus families are represented by different colors and symbols. The red solid line shows the power law of Eq. (3) for the expected values of  $\langle \text{MLD} \rangle$  for random RNA sequences, constrained only by their overall viral-like nucleotide composition. Due to their atypical nucleotide composition, Tymoviridae are not represented by Eq. (3), and the corresponding scaling law for Tymoviridae-like random RNA sequences,  $\overline{\langle \text{MLD} \rangle}_{\text{T}_y}(N) = (0.92 \pm 0.44) \times N^{0.669 \pm 0.054}$ , is shown with an orange dashed line. See SI for further information. To see this figure in color, go online.

## RESULTS

### Validation: compactness of WT and random RNA sequences

As a starting point for our analysis we considered an extensive set of 128 WT viral sequences listed in supporting information (SI) Table S1. We characterised their compactness by following the method introduced by Yoffe et al. [33], which entails two steps, detailed in the Materials and Methods section. The first step consists of computing an ensemble of several hundred representative planar RNA folds using the ViennaRNA package [49]. Next, one calculates the maximum ladder distance (MLD) of each fold. We recall that the ladder distances are obtained by considering in turn all possible pairs of nucleotides and identifying their shortest connecting path, i.e., the one with the minimal number of “rungs on the ladder” along the duplexed parts of the folds. The number of rungs of the longest minimal path is the MLD, an example of which is shown in Fig. 1(a).

As discussed in Refs. [33] and [34], the thermal average of the MLD, denoted by  $\langle \text{MLD} \rangle$ , is a viable, albeit coarse-grained proxy for the equilibrium spatial compactness of a folded sequence. Since it can be calculated by highly efficient algorithms, it is particularly apt for numerical implementation in extensive enumerative contexts such as the present one.

The comparison of the  $\langle \text{MLD} \rangle$ s computed for the 128 viral sequences considered in our study with the  $\langle \text{MLD} \rangle$ s of random sequences with viral-like nucleotide composition (see Materials and Methods) conforms to the earlier conclusion of Yoffe et al. [33] that WT RNA genomes have an enhanced fold compactness compared to arbitrary RNA sequences. This point is illustrated in Figs. 1(b) and (c). As can be seen in Fig 1(c), the  $\langle \text{MLD} \rangle$ s of random RNA sequences, additionally averaged over several possible mutations, follow the power law

$$\overline{\langle \text{MLD} \rangle}(N) = (1.365 \pm 0.05) \times N^{0.662 \pm 0.004}, \quad (3)$$

where the overline indicates the additional averaging over different possible mutations. On the other hand, the  $\langle \text{MLD} \rangle$ s of WT sequences are almost always more compact than the corresponding random values given by Eq. (3). We also note that the parameters of the power law given by Eq. (3) are in good accord with the findings of Ref. [33].

### Compactness of WT and synonymously-mutated RNA sequences

Since the fixation of mutations in viral genomes is subject to a number of evolutionary pressures, the fact that WT RNA sequences of icosahedral viruses tend to be more compact than predicted by Eq. (3) is not enough to conclude

that they have been evolutionary selected for optimal compactness. In fact, the sequence space accessible to random mutations is unrealistically large because it does not account for the several selection constraints that viable RNA sequences have to obey.

Arguably, the most severe of such constraints reflects the necessity for the viruses to preserve their protein phenotype. Accordingly, we explore its implications for genome compactness by considering only sequences which encode for the same proteins as the WT RNA. This amounts to restricting our considerations only to the rather limited combinatorial subspace of *synonymous variants* of WT viral RNA sequences.

We recall that synonymous mutations originate in the degenerate mapping of the 61 possible codons, which are nucleotide triplets, to the 20 canonical amino acids. Equivalent codons typically differ only at the third nucleotide [50]. Accordingly, we shall assume, for simplicity, that the A, U, G, and C nucleotides can appear with equal probability at the third codon position, one can estimate that two synonymous versions of a gene have a nucleotide sequence identity of about 75 %. Since, in the set of viruses considered in our study, on average ( $90 \pm 7$ ) % of the genome codes for at least one gene, and additionally assuming for simplicity that the four nucleotides have equal probability in the non-coding region which we are not constraining, we can estimate that at least around 66-73 % of the whole genome will be conserved under synonymous mutation flow.

This limited genome composition variability is further thinned down both by the imposed conservation of the dinucleotide composition characteristic for the virus family and, in some viruses, by the presence of overlapping reading frames which dramatically reduce the possibility to mutate the third nucleotide in a codon. Due to these two factors, it is found that typical sequence identity between WT sequences and their synonymous mutations ranges from about 66 % to 85 %, as shown in Fig. 2(b).

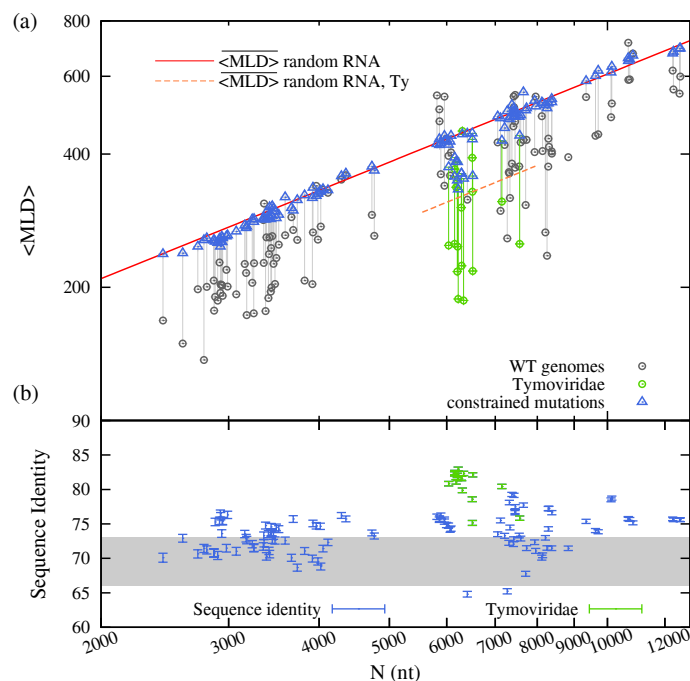


Figure 2. a) Influence of synonymous point mutations on MLD. The  $\langle \text{MLD} \rangle$ s of WT viral sequences from Fig. 1(b) are shown as gray circles, and the  $\langle \text{MLD} \rangle$ s of synonymously mutated sequences as blue triangles. Scaling laws for  $\langle \text{MLD} \rangle$  of random RNA sequences with viral-like and Tymoviridae-like composition are shown as in Fig. 1. b) The average degree of sequence identity between the mutated and WT sequences. The gray shaded area indicates the values one would expect if only one in three nucleotides were allowed to mutate in the coding regions of the genomes. Note that Tymoviridae genomes, marked in green, are more conserved than the others. This is due to the presence of overlapping reading frames covering on average 30 % of their genome. To see this figure in color, go online.

The sequence space of synonymous mutations is thus so severely restricted that there is no reason to expect that their progressive accumulation has the same effect on compactness as the unrestricted random shuffling of viral RNA sequences. As a matter of fact, the constrained synonymously-mutated sequences could have, *a priori*, about the same compactness as WT sequences or even improve it! To support the earlier observations that WT RNAs are optimized for their spatial compactness, one must therefore necessarily demonstrate that the accumulation of synonymous mutations, while leaving the encoded protein phenotype and the chemical composition of the sequence unchanged,



progressively destroys the spatial compactness observed in WT sequences and quantified by their respective MLDs.

To address this point, we start from WT viral RNA sequences and generate a mutation flow in the sequence space using a Monte Carlo (MC) algorithm which proposes point mutations of the sequence and accepts or rejects them based on the constraints of synonymity and the conservation of the dinucleotide frequencies characteristic for a given virus family (see also Materials and Methods). The typical compactness of the resulting synonymously mutated WT genomes is again characterized by the asymptotic value of  $\langle \text{MLD} \rangle$ , averaged additionally over different mutated sequences and denoted by  $\overline{\langle \text{MLD} \rangle}$ .

The resulting MLDs are shown in Fig. 2(a). It is indeed striking to notice that despite the strongly reduced available sequence space, the  $\overline{\langle \text{MLD} \rangle}$  of synonymously mutated sequences falls on the same curve which describes the  $\overline{\langle \text{MLD} \rangle}$  of random sequences, given by the power law in Eq. (3). This fundamental observation can be condensed in the symbolic statement:

$$\overline{\langle \text{MLD} \rangle}_{\text{WT}}(N) \xrightarrow{(\text{syn})} \overline{\langle \text{MLD} \rangle}_{\text{random}}(N), \quad (4)$$

where  $N$  is the genome length and the arrow is a shorthand for indicating the flow in the synonymous mutations subspace.

This result proves the conjecture that the WT genomes are indeed characterized by a certain optimality of the MLD which, in turn, reflects atypically-high degrees of RNA fold compactness. In fact, the results of Fig. 2(b) demonstrate that the WT MLD/compactness can be obliterated even within a much restricted subset of mutations that otherwise leave the viral phenotype and sequence composition unchanged.

As an aside, we note that Tymoviridae exhibit an atypical behavior, with the limiting value of  $\overline{\langle \text{MLD} \rangle}$  under the synonymous mutation flow approaching values which are still below the ones characteristic for random RNAs. The reason for this lies in the fact that Tymoviridae have a different nucleotide composition with respect to other viral families; accounting for this different composition one obtains a different prefactor for the scaling law in Eq. (3), corresponding to more compact values of MLD, as shown in Fig. 1(c); see also SI Fig. S3 for more details.

### Synonymous mutation flow and the stability of genome MLD

The previous result leads us to examine the details of the implied synonymous mutation flow [Eq. (4)] and the stability of the terminal, asymptotic state of the mutated sequence. In particular, we wish to establish the minimal number of point nucleotide mutations that are needed to bring the MLD of a viral RNA from its WT value to the random reference value. It is especially interesting to ascertain whether this change in compactness happens progressively, indicating that a continuous accumulation of mutations is responsible for disrupting the WT RNA spatial compactness, or whether the change is due to sporadic, punctuated events, which would suggest the presence of specific RNA “hotspots” where mutations can dramatically affect fold compactness.

To illuminate this point we considered 9 synthetic synonymous mutation flow trajectories for 4 different viral sequences extracted from 3 viruses picked at random from 3 different families: brome mosaic virus (BMV), ononis yellow mosaic virus (OnYMV), and equine rhinitis B virus 1 (ERBV1). The considered sequences were chosen in order to probe the whole range of genome lengths spanning from  $N \simeq 2800$  nt to  $N \simeq 8800$  nt. The trajectories were generated using the same MC scheme used to generate the equilibrium data presented in Fig. 2 (see also Materials and Methods), but with a much more frequent sampling of the mutated sequences (every  $N/100$  attempted synonymous mutations) so as to leave detectable correlations in the series of generated sequences – in this way mimicking the viral mutation dynamics.

The results are shown in Fig. 3. From the mutation flow trajectories we discern that, at least for the sequences considered, the change in compactness follows the continuous and gradual accumulation of synonymous mutations, and does not take place in a punctuated manner. Nonetheless, not many mutations are needed to make the MLD of these sequences already indistinguishable from that of randomized RNAs. In fact, mutating not more than  $\sim 5\%$  of the full genome suffices to erase the characteristic WT RNA compactness imprint.

A further interesting point clarified by the mutation flow trajectories shown in Fig. 3 is that the genome fold compactness is not completely optimized even in the case of WT sequences. In fact, for all the 4 sequences considered in Fig. 3 one occasionally observes more compact folded states, particularly during the initial part of the trajectories.

To better explore this interesting observation, we computed the probability density of finding mutated sequences with given  $\langle \text{MLD} \rangle$  as a function of the sequence identity to the WT sequence ratio, and plotted it as a color-coded heatmap. These probability density plots are shown in Fig. 4, and we can observe that for some of the genomes considered, such as BMV RNA1 and ERBV1, more compact structures are reachable even when almost all the unconstrained nucleotides have already been mutated. This point is most relevant in the present context. In fact,

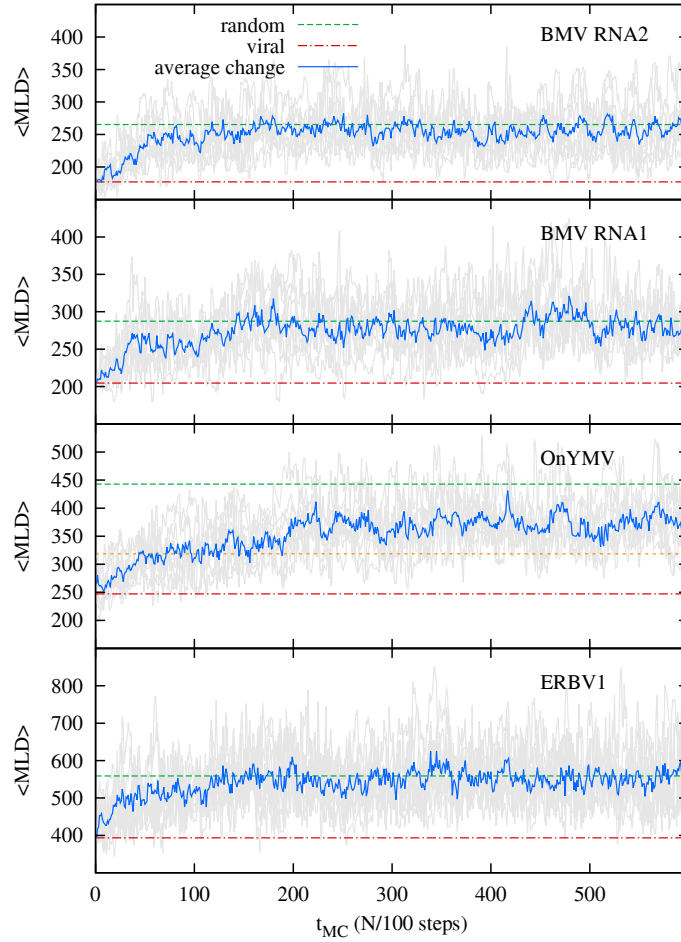


Figure 3. Mutation dynamics trajectories for 4 viral ssRNA sequences. From top to bottom: BMV RNA2 and RNA1 segments from the tripartite genome of BMV (Bromoviridae), OnYMV (Tymoviridae), and ERBV1 (Picornaviridae). Each panel shows 9  $\langle \text{MLD} \rangle$  trajectories and their average value (blue) for each sequence in units of MC steps,  $N/100$ . Red dot-dashed lines and green dashed lines show respectively the  $\langle \text{MLD} \rangle$  values of WT RNAs and the  $\langle \text{MLD} \rangle$  values of random RNAs [for viral-like composition, Eq. (3)]. Notice that in the case of OnYMV, a Tymovirus, we must consider the appropriate asymptotic value of  $\langle \text{MLD} \rangle$  for random RNAs with Tymoviridae-like composition (see Fig. 1). This value is shown by the orange short-dashed line. To see this figure in color, go online.

it demonstrates that the sequence-based synonymity constraint and the structure-based one for fold compactness, despite being in competition, can still be compatible.

This point is made more poignantly by considering the near-native pool of synonymous sequences (e.g., those with sequence identity  $\geq 95\%$ ) for the four cases presented in Fig. 4. Across these instances it is found that from 12 % to 21 % of the near-native synonymous sequences have a predicted fold compactness that is equal or higher than the wild-type one. This indicates that the well-optimised viral sequences still have a portion of phase space available for evolving while respecting both sequence- and structure-based stringent constraints. This appreciable residual mutation freedom may be clearly necessary to simultaneously accommodate other concurrent selection constraints.

#### Taking into account codon usage bias and untranslated regions

Finally, we examine the effect of two additional constraints which are known to be relevant for some viruses, and may play a role in maintaining viral RNA compactness. The first constraint is given by the presence of functionally important secondary RNA structures in the untranslated regions (UTRs) at the 3' and 5' ends of several viral



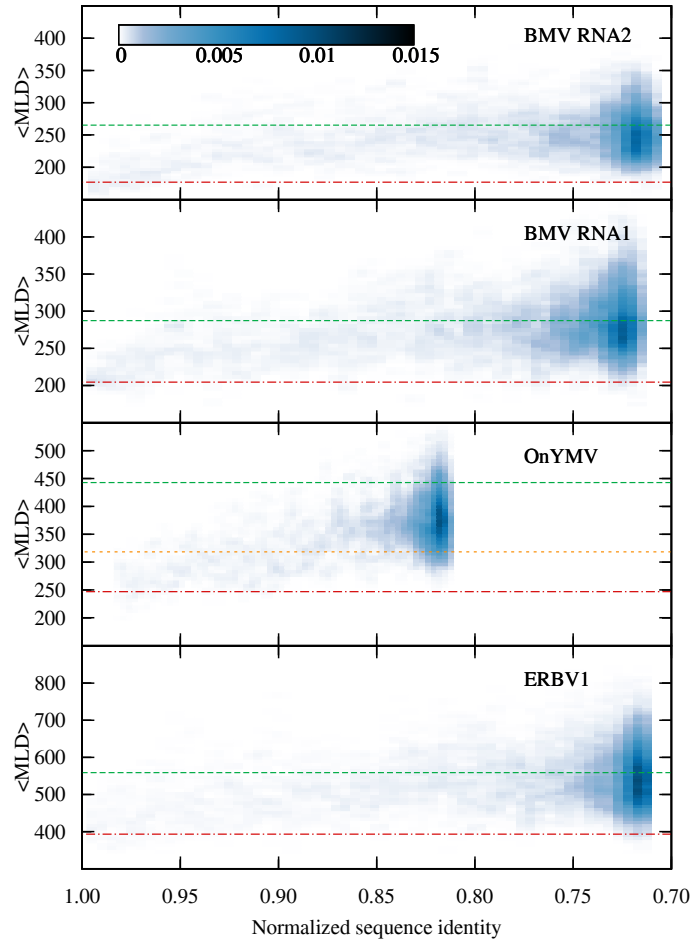


Figure 4. Color-coded heat maps for the probability density of finding mutated sequences with given  $\langle \text{MLD} \rangle$  and sequence identity with the WT sequence. The probability density for each virus is computed and normalized over the whole length of the 9 mutation trajectories (1500 MC steps) shown in Fig. 3. Red dot-dashed lines and green dashed lines show respectively the  $\langle \text{MLD} \rangle$  values of WT RNA and the  $\langle \text{MLD} \rangle$  values of random RNAs [with viral-like composition, Eq. (3)]. The orange short-dashed line in the OnYMV case shows the random  $\langle \text{MLD} \rangle$  value for Tymoviridae-like composition. To see this figure in color, go online.

genomes [51–53]. We take into account this constraint by simply limiting the mutation flow to the coding regions of the genomes. Note that with this additional constraint our theoretical estimate of the overall sequence identity between WT sequences and sequences mutated asymptotically to saturation moves from 66-73 % to 76-83 %.

The second additional constraint is given by the fact that, since viruses adapt to their hosts, not all the codons which translate into the same amino-acid are statistically equivalent: some of them are more probable than others. This codon usage bias is known to be an important constraint for several viruses. In fact, changing the codon bias or the codon-pair bias leads to attenuated viruses and has been proposed as a possible vaccination strategy [54, 55]. To produce mutated sequences with WT codon populations we shuffled the equivalent codons within every viral gene (see Materials and Methods for details regarding the implementation of codon-bias preserving synonymous mutations).

The results obtained with both of these constraints are compared in Fig. 5 against those previously obtained using synonymous point mutations. It is important to notice that even with these additional constraints, which further thin out the phase space available to mutations, our results remain valid, confirming the presence of an evolutionary pressure to produce compact RNA folds.

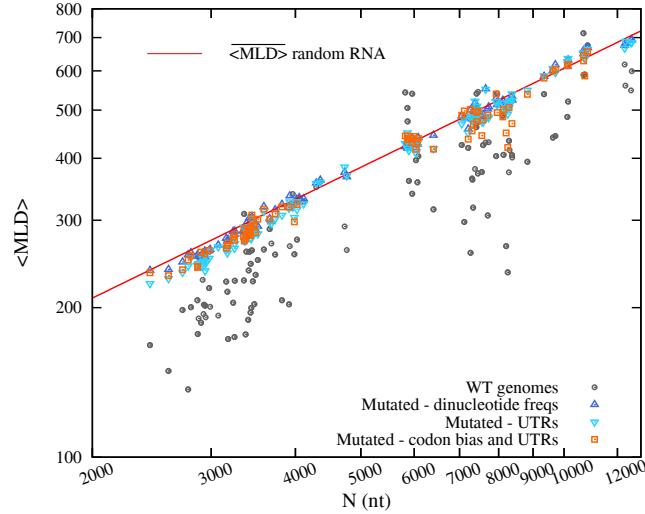


Figure 5. The  $\langle \text{MLD} \rangle$  values for the synonymous constraint only (upward triangles) and for the additional constraints of preserving UTRs sequences (downward triangles) and UTRs sequences as well as codon biases (squares). The  $\langle \text{MLD} \rangle$  values for these last two cases are evaluated over a set of 150 mutated sequences for each virus. Data are presented in the same manner as in Fig. 2 (see also SI Fig. S3 for UTRs preserving synonymous point mutations of Tymoviridae). To see this figure in color, go online.

## DISCUSSION AND CONCLUSIONS

While the fundamental mechanisms by which point mutations affect the fitness of the organisms in their respective environments (*via* the transcription of the mutated nucleotide sequence into the modified protein products) are well understood [12–14], it is less known what are their effects on the purely physico-chemical properties of their genomes. In order to investigate possible parallel selection mechanisms and eventual embedded levels of coding that control the compactness of viral ssRNA folds, we analyzed a synthetic model for accumulating synonymous mutations in viral RNAs and assessed their impact on the spatial compactness of the genome as quantified by the MLD measure, introduced by Yoffe and coworkers [33]. We have analyzed the effects of synonymous mutations under different constraints on ssRNA genomes for a large number of different viral families with icosahedral capsids, and compared the changes in their compactness with randomly shuffled RNA sequences with the same nucleotide composition, which are in general significantly less compact than those encapsidated by viruses.

Using extensive computational analysis we have shown that progressive accumulation of synonymous point mutations, although neutral from the functional point of view as they conserve the expressed protein complement, completely erases the typical compactness of viral WT RNA folds. In fact, under the synonymous mutation flow the MLDs of WT RNAs approach their corresponding random RNA values in a continuous manner even after a relatively small number of mutations. Although, in principle, the emergence of viral RNA fold compactness may still be related to some other evolutionary pressure, our results rule out the principal ones, including codon bias and the preservation of functional UTRs, and thus strongly support the independent evolution of viral RNA fold compactness. Arguably, such a dramatic reduction in RNA fold compactness, which in this respect eventually makes it undistinguishable from a random RNA sequence, has a relevant impact on the virion assembly and therefore on the ability of viruses to replicate and propagate their infection. These results are strengthened by the observation that the typical WT RNA compactness is not related to codon usage bias nor is it dictated by the particular sequence/structure of its non-coding regions, since synonymous mutations which maintain both unchanged were found to nonetheless destroy the typical WT RNA compactness.

The connection between the viral RNA sequence and its physical properties, such as its compactness, may in future allow to control the physical properties of viral RNAs and specifically their aptitude for efficient packing. This we believe may lead to improve and broaden the scope of existing strategies which harness viral mutation rates to achieve virus attenuation.

## SUPPORTING MATERIAL

Seven figures and two tables of supporting data.

LT, ALB, and RP acknowledge support from ARRS Grants No. P1-0055 and J1-4297. CM acknowledges support from the Italian Ministry of Education, grant PRIN No. 2010HXAW77.

- 
- [1] Olsthoorn, R. C., and J. Van Duin. 1996. Evolutionary reconstruction of a hairpin deleted from the genome of an RNA virus. *Proc. Natl. Acad. Sci. USA* 93:12256–12261.
  - [2] Klovins, J., V. Berzins, and J. Van Duin. 1998. A long-range interaction in Qbeta RNA that bridges the thousand nucleotides between the M-site and the 3' end is required for replication. *RNA* 4:948–957.
  - [3] Dykeman, E. C., P. G. Stockley, and R. Twarock. 2013. Packaging signals in two single-stranded RNA viruses imply a conserved assembly mechanism and geometry of the packaged genome. *J. Mol. Biol.* 425:3235–3249.
  - [4] Dykeman, E. C., P. G. Stockley, and R. Twarock. 2014. Solving a Levinthal's paradox for virus assembly identifies a unique antiviral strategy. *Proc. Natl. Acad. Sci. USA* doi:10.1073/pnas.1319479111.
  - [5] Simmonds, P., A. Tuplin, and D. J. Evans. 2004. Detection of genome-scale ordered RNA structure (GORS) in genomes of positive-stranded RNA viruses: Implications for virus evolution and host persistence. *RNA* 10:1337–1351.
  - [6] Sanjuán, R., and A. V. Bordería. 2011. Interplay between RNA structure and protein evolution in HIV-1. *Mol. Biol. Evol.* 28:1333–1338.
  - [7] Cuevas, J. M., P. Domingo-Calap, and R. Sanjuán. 2012. The fitness effects of synonymous mutations in DNA and RNA viruses. *Mol. Biol. Evol.* 29:17–20.
  - [8] Davis, M., S. M. Sagan, J. P. Pezacki, D. J. Evans, and P. Simmonds. 2008. Bioinformatic and physical characterization of genome-scale ordered RNA structure in mammalian RNA viruses. *J. Virol.* 82:11824–11836.
  - [9] Holmes, E. C. 2009. *The Evolution and Emergence of RNA Viruses*. Oxford University Press, New York.
  - [10] Belshaw, R., A. Gardner, A. Rambaut, and O. G. Pybus. 2007. Pacing a small cage: mutation and RNA viruses. *Trends Ecol. Evol.* 23:188–193.
  - [11] Eigen, M. 2000. Viruses: Evolution, propagation, and defense. *Nutr. Rev.* 58(s1):S5–S16.
  - [12] Wylie, C. D., and E. I. Shakhnovich. 2011. A biophysical protein folding model accounts for most mutational fitness effects in viruses. *Proc. Natl. Acad. Sci. USA* 108:9916–9921.
  - [13] Chen, P., and E. I. Shakhnovich. 2009. Lethal mutagenesis in viruses and bacteria. *Genetics* 183:639–650.
  - [14] Duffy, S., L. A. Shackleton, and E. C. Holmes. 2008. Rates of evolutionary change in viruses: patterns and determinants. *Nat. Rev. Genet.* 9:267–276.
  - [15] Gong, L. I., M. A. Suchard, and J. D. Bloom. 2013. Stability-mediated epistasis constrains the evolution of an influenza protein. *eLife* 2:e00631.
  - [16] Hyeon, C., R. I. Dima, and D. Thirumalai. Size, shape and flexibility of RNA structures. *J. Chem. Phys.* 125:194905.
  - [17] Marenduzzo, D., E. Orlandini, A. Stasiak, D. W. Sumners, L. Tubiana, and C. Micheletti. 2009. DNA-DNA interactions in bacteriophage capsids are responsible for the observed DNA knotting. *Proc. Natl. Acad. Sci. USA* 106:22269.
  - [18] Marenduzzo, D., C. Micheletti, E. Orlandini, and D. W. Sumners. 2013. Topological friction strongly affects viral DNA ejection. *Proc. Natl. Acad. Sci. USA* 110:20081–20086.
  - [19] Erdemci-Tandogan, G., J. Wagner, P. van der Schoot, R. Podgornik, and R. Zandi. 2014. RNA topology remodels electrostatic stabilization of viruses. *Phys. Rev. E* 89:032707.
  - [20] Nap, R. J., A. Lošdorfer Božič, I. Szleifer, and R. Podgornik. 2014. The role of solution conditions in the bacteriophage PP7 capsid charge regulation. *Biophys. J.* in print.
  - [21] Caspar, D. L. D., and K. Namba. 1990. Switching in the self-assembly of tobacco mosaic virus. *Adv. Biophys.* 26:157–185.
  - [22] Bruinsma, R. F., W. M. Gelbart, D. Reguera, J. Rudnick, and R. Zandi. 2003. Viral self-assembly as a thermodynamic process. *Phys. Rev. Lett.* 90:248101.
  - [23] Reguera, J., A. Carreira, L. Riobobos, J. M. Almendral, and M. G. Mateu. 2004. Role of interfacial amino acid residues in assembly, stability, and conformation of a spherical virus capsid. *Proc. Natl. Acad. Sci. USA* 101:2724–2729.
  - [24] Singh, S., and A. Zlotnick. 2003. Observed hysteresis of virus capsid disassembly is implicit in kinetic models of assembly. *J. Biol. Chem.* 278:18249–18255.
  - [25] Nguyen, H. D., V. S. Reddy, and C. L. Brooks III. 2007. Deciphering the kinetic mechanism of spontaneous self-assembly of icosahedral capsids. *Nano. Lett.* 7:338–344.
  - [26] Castellanos M., R. Pérez, C. Carrasco, M. Hernando-Pérez, J. Gómez-Herrero, P. J. de Pablo, and M. G. Mateu. 2012. Mechanical elasticity as a physical signature of conformational dynamics in a virus particle. *Proc. Natl. Acad. Sci. USA* 109:12028–12033.
  - [27] Roos, W. H., I. Gertsman, E. R. May, C. L. Brooks III, J. E. Johnson, and G. J. L. Wuite. 2012. Mechanics of bacteriophage maturation. *Proc. Natl. Acad. Sci. USA* 109:2342–2347.
  - [28] Polles, G., G. Indelicato, R. Potestio, P. Cermelli, R. Twarock, and C. Micheletti. 2013. Mechanical and assembly units of viral capsids identified via quasi-rigid domain decomposition. *PLoS Comp. Bio.* 9:e1003331.
  - [29] Cadena-Nava, R. D., M. Comas-Garcia, R. F. Garmann, A. L. N. Rao, C. M. Knobler, and W. M. Gelbart. 2012. Self-assembly of viral capsid protein and RNA molecules of different sizes: requirement for a specific high protein/RNA mass

- ratio. *J. Virol.* 86:3318–3326.
- [30] Comas-Garcia, M., R. D. Cadena-Nava, A. L. N. Rao, C. M. Knobler, and W. M. Gelbart. 2012. In vitro quantification of the relative packaging efficiencies of single-stranded RNA molecules by viral capsid protein. *J. Virol.* 86:12271–12282.
  - [31] Perlmutter, J. D., Q. Cong, and M. F. Hagan. 2013. Viral genome structures are optimal for capsid assembly. *eLife* 2:e00632.
  - [32] Harvey, S. C., Y. Zeng, and C. E. Heitsch. 2013. The icosahedral RNA virus as a grotto: organizing the genomes into stalagmites and stalactites. *J. Biol. Phys.* 39:163–172.
  - [33] Yoffe, A. M., P. Prinsen, A. Gopal, C. M. Knobler, W. M. Gelbart, and A. Ben-Shaul. 2008. Predicting the sizes of large RNA molecules. *Proc. Natl. Acad. Sci. USA* 105:16153.
  - [34] Fang, L. T., W. M. Gelbart, and A. Ben-Shaul. 2011. The size of RNA as an ideal branched polymer. *J. Chem. Phys.* 135:155105.
  - [35] Gopal, A., Z. H. Zhou, C. M. Knobler, and W. M. Gelbart. 2012. Visualizing large RNA molecules in solution. *RNA* 18:284–299.
  - [36] Biotechnology, National Center for. *NCBI Nucleotide Database*. June–July 2010. Web. 08 July 2010. <http://www.ncbi.nlm.nih.gov/nucleotide>
  - [37] Hulo, C., E. De Castro, P. Masson, L. Bougueleret, A. Bairoch, I. Xenarios, and P. Le Mercier. 2011. ViralZone: a knowledge resource to understand virus diversity. *Nucleic Acids Res.* 39:D576–D582.
  - [38] Simon-Loriere, E., E. C. Holmes, and I. Págan. 2013. The effect of gene overlapping on the rate of RNA virus evolution. *Mol. Biol. Evol.* 30:1916–1928.
  - [39] Chirico, N., A. Vianelli, and R. Belshaw. 2010. Why genes overlap in viruses. *Proc. R. Soc. B* 277:3809–3817.
  - [40] Pedersen, A.-M. K., and J. L. Jensen. 2001. A dependent-rates model and an MCMC-based methodology for the maximum-likelihood analysis of sequences with overlapping reading frames. *Mol. Biol. Evol.* 18:763–776.
  - [41] Chung, W.-Y., S. Wadhawan, R. Szklarczyk, S. Kosakovsky Pond, and A. Nekrutenko. 2007. A first look at ARFome: dual-coding genes in mammalian genomes. *PLoS Comput. Biol.* 3:e91.
  - [42] Nussinov, R. 1981. Nearest neighbor nucleotide patterns. Structural and biological implications. *J. Biol. Chem.* 256:8458–8462.
  - [43] Durstenfeld, R. 1964. Algorithm 235: random permutation. *Commun. ACM* 7:420.
  - [44] Knuth, D. E. 1981. Seminumerical Algorithms. Vol. 2 of The Art of Computer Programming. Addison-Wesley, Reading, Massachusetts, 2nd ed.
  - [45] Anisimova, M., and C. Kosiol. 2009. Investigating protein-coding sequence evolution with probabilistic codon substitution models. *Mol. Biol. Evol.* 26:255–271.
  - [46] Gu, W., T. Zhou, and C. O. Wilke. 2010. A universal trend of reduced mRNA stability near the translation-initiation site in prokaryotes and eukaryotes. *PLoS Comput. Biol.* 6:e1000664.
  - [47] Saito, M., and M. Matsumoto. 2008. SIMD-oriented fast Mersenne Twister: a 128-bit pseudorandom number generator. In Monte Carlo and Quasi-Monte Carlo Methods 2006. A. Keller, S. Heinrich, and H. Niederreiter, editors. Springer, Berlin Heidelberg. 607–622.
  - [48] Bundschuh, R., and T. Hwa. 2002. Statistical mechanics of secondary structures formed by random RNA sequences. *Phys. Rev. E* 65:031903.
  - [49] Lorenz, R., S. H. Bernhart, C. Höner zu Siederdissen, H. Tafer, C. Flamm, P. F. Stadler, and I. L. Hofacker. 2011. ViennaRNA Package 2.0. *Algorithm. Mol. Biol.* 6:26.
  - [50] Nelson, D. L., and M. M. Cox. 2008. Lehninger Principles of Biochemistry. W. H. Freeman and Company, New York, 5th ed.
  - [51] Marz, M., N. Beerenwinkel, C. Drosten, M. Fricke, D. Frishman, I. L. Hofacker, D. Hoffmann, M. Middendorf, T. Rattei, P. F. Stadler, and A. Töpfer. 2014. Challenges in RNA virus bioinformatics. *Bioinformatics* 30:1793–1799.
  - [52] Alvarez, D. E., A. L. De Lella Ezcurra, S. Fucito, and A. V. Gamarnik. 2005. Role of RNA structures present at the 3' UTR of dengue virus on translation, RNA synthesis, and viral replication. *Virology* 399:200–212.
  - [53] Tsukiyama-Kohara, K., N. Iizuka, M. Kohara, and A. Nomoto. 1992. Internal ribosome entry site within hepatitis C virus RNA. *J. Virol.* 66:1476–1483.
  - [54] Bull, J. J., I. J. Molineux, and C. O. Wilke. 2012. Slow fitness recovery in a codon-modified viral genome. *Mol. Biol. Evol.* 24:1–8.
  - [55] Coleman, J. R., D. Papamichail, S. Skiena, B. Fitcher, E. Wimmer, and S. Mueller. 2008. Virus attenuation by genome-scale changes in codon pair bias. *Science* 320:1784–1787.

# Supporting Information

## I. FIT OF THE SHUFFLED RNA MLD

To obtain the power law for the MLD of random RNAs, we shuffled 12 RNA sequences of different lengths (1000 nt, 1500 nt, ..., 6000 nt), all having a viral-like nucleotide composition: 0.26 A, 0.28 U, 0.24 G, 0.22 C (obtained excluding Tymoviridae, which have a significantly different composition). For every sequence length, we produced 500 independent sequences over which we computed the expected (thermally averaged)  $\langle \text{MLD} \rangle$ . The power law of Eq. (1) in the main text is then obtained by fitting the dependence of  $\langle \text{MLD} \rangle$ , further averaged over the 500 different mutations, on the sequence length.

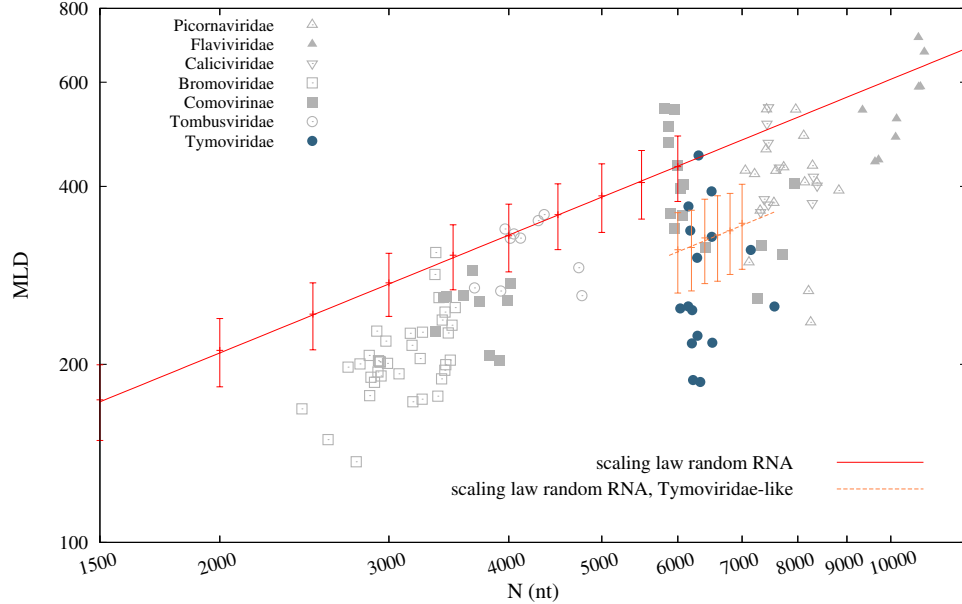
As already mentioned in the main text, Tymoviridae differ notably from the other families in their nucleotide composition, and they were not considered when producing the averaged viral-like composition. Evaluating the average composition for the set of Tymoviridae viruses considered in the main text, we obtain 0.20 A, 0.24 U, 0.18 G, 0.38 C.

Using this alternative composition and adopting the same procedure used for the other families we obtain a scaling law describing the  $\langle \text{MLD} \rangle$  dependence of Tymoviridae-like random RNA sequences:

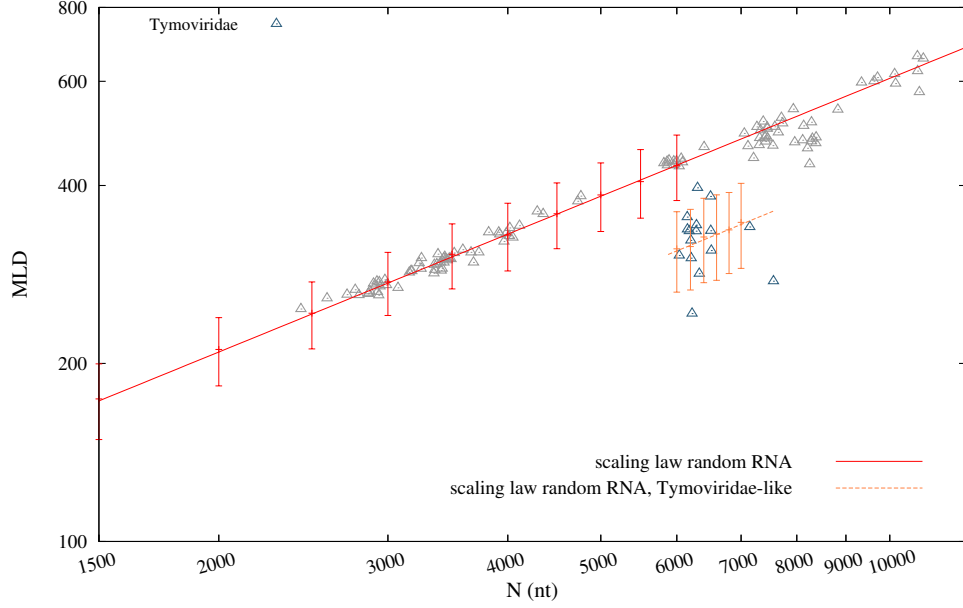
$$\overline{\langle \text{MLD} \rangle}_{\text{Ty}}(N) = (0.92 \pm 0.44) \times N^{(0.669 \pm 0.054)}. \quad (5)$$

Note that the exponent,  $0.669 \pm 0.054$ , is compatible with the one obtained for the other viral families,  $0.662 \pm 0.004$ . Both fits are shown in Fig. S1.

We further check the validity of the scaling laws for our viral families by randomly shuffling the WT RNA sequences themselves, without any further constraints. The results, shown in Fig. S2, show once again that the two scaling laws are a good reference for random RNAs with the viral-like composition considered in our sample. For Tymoviridae, we notice that a couple of viruses remain more compact than predicted by Eq. (S5). This is due to them having a composition which is substantially different from the Tymoviridae average composition.



SI Fig. 1.  $\langle \text{MLD} \rangle$  values of WT RNA genomes are shown in gray for all families apart from Tymoviridae, which are highlighted in blue.  $\langle \text{MLD} \rangle$  values of random sequences are shown with red and orange errorbars for viral-like and Tymoviridae-like nucleotide composition, respectively. The respective fitting lines are displayed with the same colors. The  $p$ -value of the fit parameters for the viral-like composition is below  $10^{-10}$ , and the adjusted  $R^2$  is 0.999948. For Tymoviridae-like composition the  $p$ -value of exponent is  $\simeq 10^{-4}$  and the adjusted  $R^2$  is 0.999968.



SI Fig. 2.  $\langle \text{MLD} \rangle$  values of randomly shuffled WT RNA genomes, shown in gray for all families apart from Tymoviridae, which are highlighted in blue.  $\langle \text{MLD} \rangle$  values of random sequences are shown with red and orange errorbars for viral-like and Tymoviridae-like nucleotide composition, respectively.

## II. MUTATIONS PRESERVING UTRS

As detailed in the main text, we further tested the robustness of our results by adding additional optional constraint as the preservation of Untranslated regions (UTRs) near the ends of the genome and the preservation of the codon biases within each gene. The  $\langle \text{MLD} \rangle$  values obtained with these additional constraint are compared with those obtained under the constraint of synonymous mutations only in Fig. 5 in the main text. Here, in Fig. S3 we extend the comparison for the additional constraint of preserving UTRs to include Tymoviridae.  $\langle \text{MLD} \rangle$  values under the additional constraint of fixed codon bias were not calculated for this family since all the tymoviridae genomes in our set present overlapping genes.

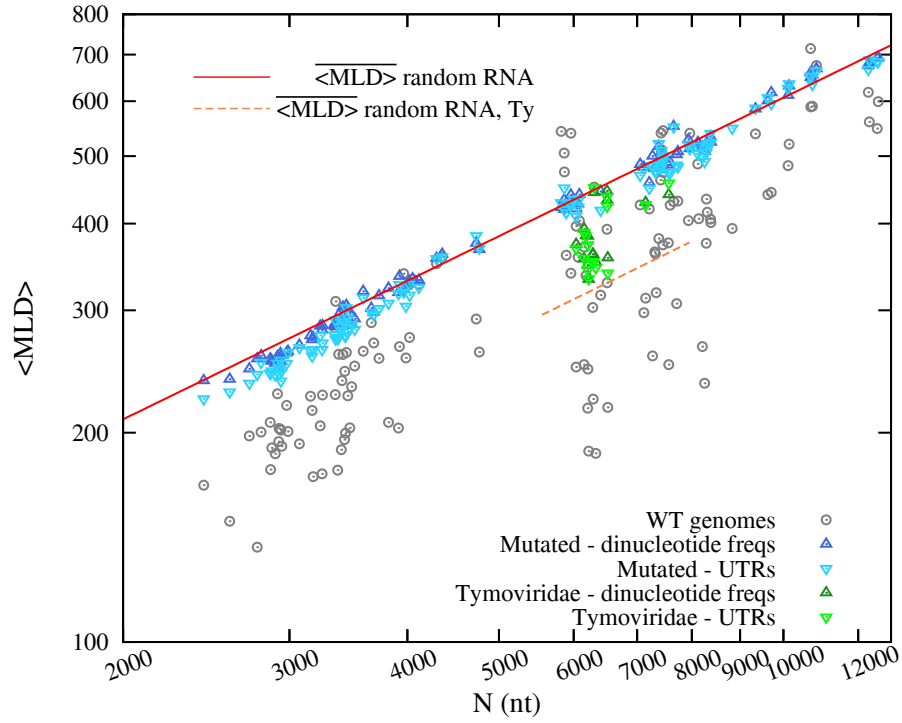
## III. MUTATIONS AT FIXED NUCLEOTIDE COMPOSITION

To test the robustness of the results reported in the main text, we implemented another mutation flow which conserves the nucleotide composition instead of the dinucleotide frequencies. This is achieved by using a Fisher-Yates algorithm where proposed shuffles are accepted or rejected on the basis of whether or not the resulting genome still encodes for the same proteins. The results of this different simulation setup are shown in Fig. S4.

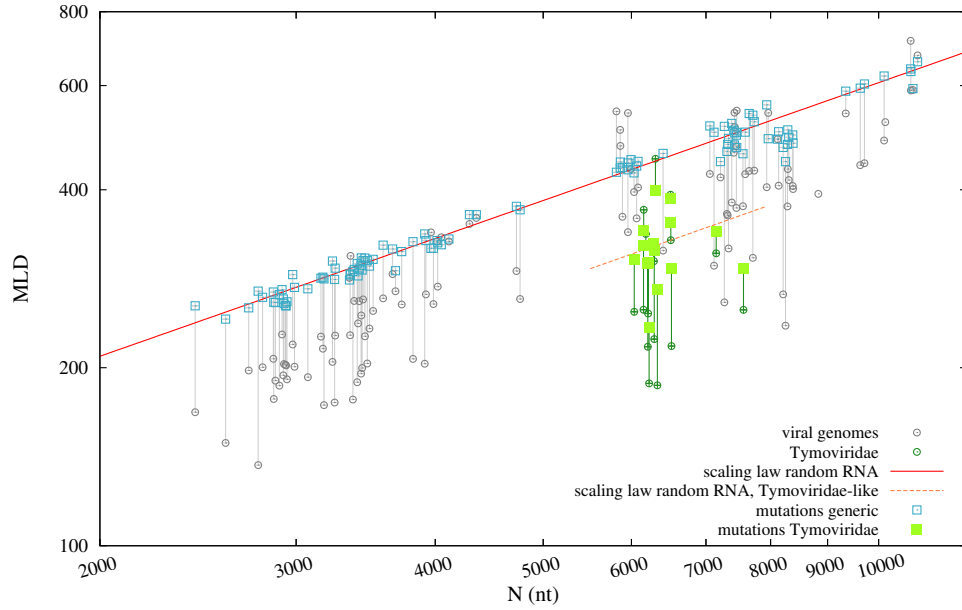
Note that the values of  $\langle \text{MLD} \rangle$  obtained in this way show a clear correlation with those obtained by unrestricted random shuffling of the WT RNA sequences, shown in Fig. S2.

## IV. DETAILS OF DINUCLEOTIDE AND NUCLEOTIDE COMPOSITIONS

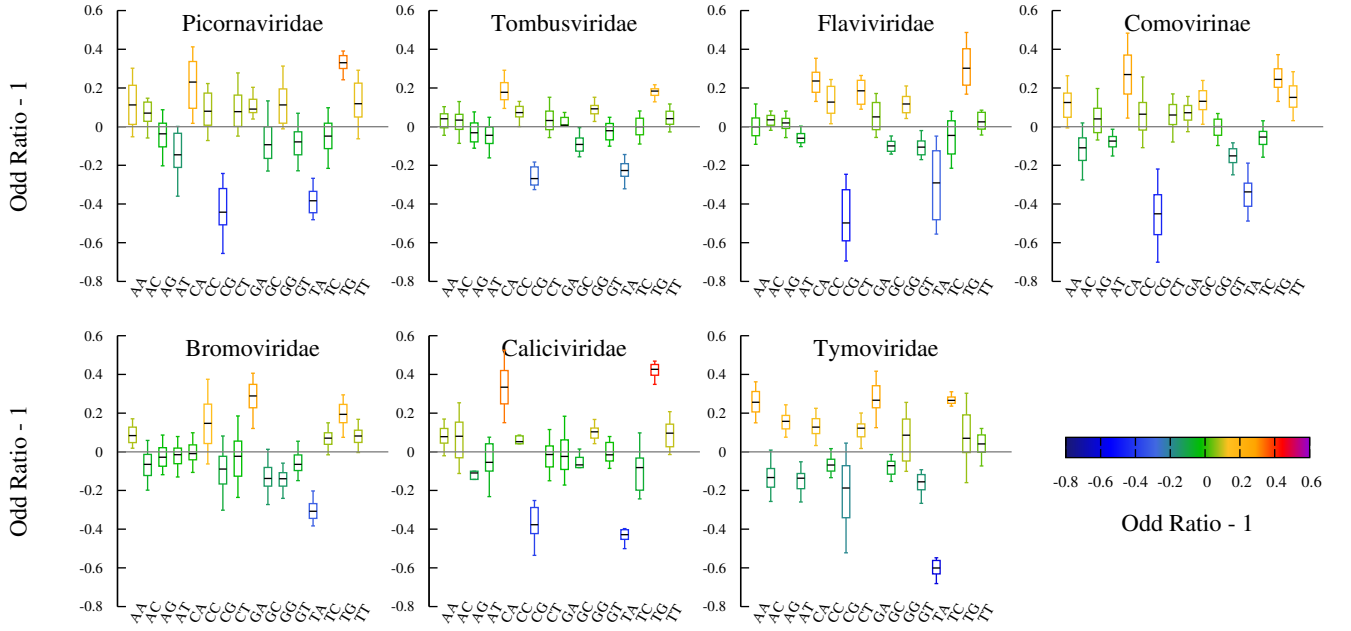




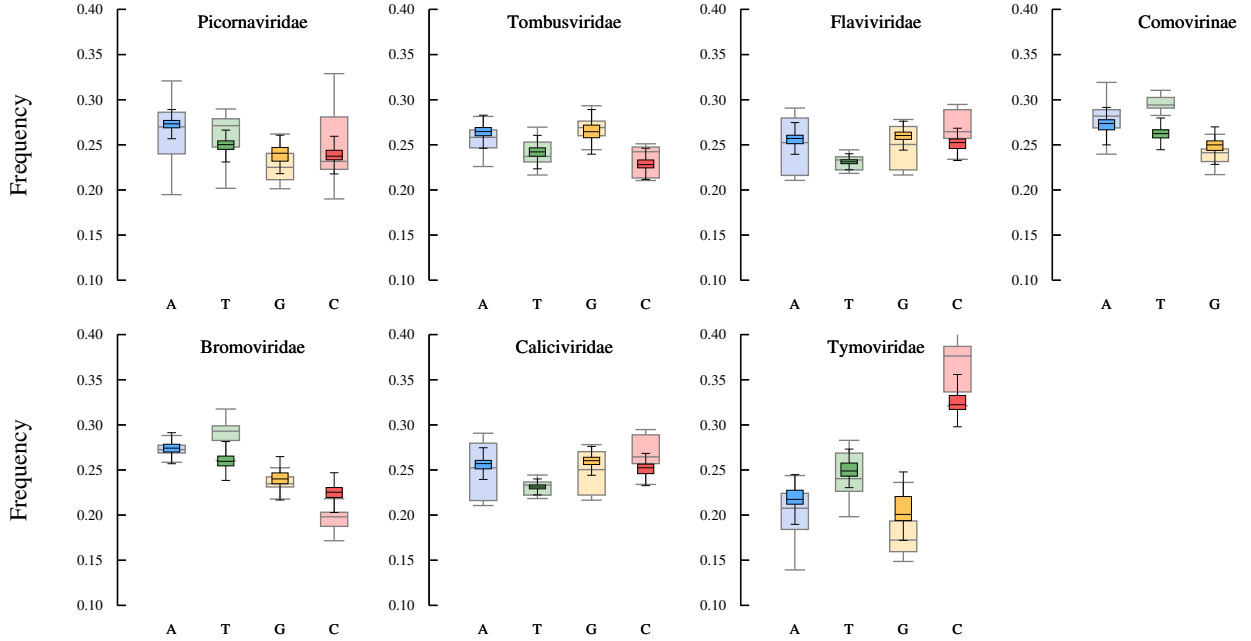
SI Fig. 3. Comparison between the  $\langle \text{MLD} \rangle$  values for the synonymous constraint only (upward triangles) and for the additional constraints of preserving UTRs sequences (downward triangles), including Tymoviridae. In the latter case  $\langle \text{MLD} \rangle$  have been evaluated over a set of 150 mutated sequences per virus.



SI Fig. 4. Mutations performed at fixed nucleotide composition. Note that the  $\langle \text{MLD} \rangle$  of the mutated viral sequences approaches the random RNA values for viral-like and Tymoviridae-like nucleotide composition in both respective cases. We note that for Tymoviridae there are some viruses which remain more compact than the corresponding random RNAs. We argue that this is due to the fact that Tymoviridae show notable fluctuations in their nucleotide composition.



SI Fig. 5. Dinucleotide odd-ratios, rescaled to zero, for the viral families considered in our study. Boxes represent quartiles, and whiskers correspond to 1.5 of the interquartile distance. These values have been used to constrain the mutation flow and produce sequences with viral-like dinucleotide frequencies, see Materials and Methods section in the main text.



SI Fig. 6. Nucleotide frequencies for WT sequences (large boxes) and mutated sequences with constrained dinucleotide composition (small boxes), shown for each virus family considered in our study. Note that in most families the imposition of conserved nucleotide frequencies results in conserved nucleotide frequencies as well, although for Tymoviridae, Comovirinae, and Bromoviridae the frequencies are not so well preserved, showing the effects of transversion changes.

## V. DATASET OF VIRAL GENOMES

Taxon	Family	Nucline code	PDB code	length	$\langle \text{MLD} \rangle_{WT}$	$\langle \text{MLD} \rangle_{mut}$	$\langle \text{MLD} \rangle_{UTRs}$	$\langle \text{MLD} \rangle_{CB}$
Bromoviridae	Anulavirus	PeZSV_RNA1	—	3383	259 ± 19	295 ± 48	274 ± 45	280 ± 43
Bromoviridae	Anulavirus	PeZSV_RNA2	—	2435	168 ± 8	238 ± 39	224 ± 37	235 ± 38
Bromoviridae	Bromovirus	BMV_RNA1	1js9	3234	204 ± 15	285 ± 46	276 ± 43	278 ± 46
Bromoviridae	Bromovirus	BMV_RNA2	1js9	2865	177 ± 12	255 ± 42	244 ± 38	242 ± 40
Bromoviridae	Bromovirus	BrBMV_RNA1	—	3158	225 ± 21	276 ± 44	263 ± 43	264 ± 42
Bromoviridae	Bromovirus	BrBMV_RNA2	—	2799	200 ± 17	258 ± 43	252 ± 40	255 ± 43
Bromoviridae	Bromovirus	CaYBV_RNA1	—	3178	172 ± 14	274 ± 45	263 ± 40	270 ± 45
Bromoviridae	Bromovirus	CaYBV_RNA2	—	2720	197 ± 20	247 ± 41	236 ± 40	239 ± 39
Bromoviridae	Bromovirus	CCMV_RNA1	1cwp	3171	215 ± 27	272 ± 44	258 ± 40	267 ± 44
Bromoviridae	Bromovirus	CCMV_RNA2	1cwp	2774	136 ± 11	256 ± 43	243 ± 38	249 ± 40
Bromoviridae	Bromovirus	MeYFV_RNA1	—	3249	174 ± 19	281 ± 46	263 ± 44	274 ± 47
Bromoviridae	Bromovirus	MeYFV_RNA2	—	2862	207 ± 11	255 ± 40	243 ± 38	241 ± 37
Bromoviridae	Bromovirus	SpBLV_RNA1	—	3252	226 ± 25	285 ± 45	269 ± 43	281 ± 47
Bromoviridae	Bromovirus	SpBLV_RNA2	—	2898	186 ± 16	252 ± 41	242 ± 40	259 ± 39
Bromoviridae	Cucumovirus	GaMMV_RNA1	—	3350	283 ± 28	286 ± 47	277 ± 43	278 ± 44
Bromoviridae	Cucumovirus	GaMMV_RNA2	—	2935	202 ± 8	260 ± 43	254 ± 39	257 ± 42
Bromoviridae	Cucumovirus	PeSV_RNA1	—	3357	309 ± 16	287 ± 46	273 ± 44	273 ± 42
Bromoviridae	Cucumovirus	TAV_RNA1	1laj	3410	237 ± 18	287 ± 46	285 ± 44	281 ± 49
Bromoviridae	Cucumovirus	TAV_RNA2	1laj	3074	192 ± 18	267 ± 43	265 ± 42	—
Bromoviridae	Ilarvirus	ApMV_RNA1	—	3476	203 ± 36	297 ± 49	283 ± 45	292 ± 49
Bromoviridae	Ilarvirus	ApMV_RNA2	—	2979	218 ± 20	261 ± 43	251 ± 43	261 ± 45
Bromoviridae	Ilarvirus	CiLRV_RNA1	—	3404	189 ± 27	289 ± 46	281 ± 48	289 ± 4
Bromoviridae	Ilarvirus	CiLRV_RNA2	—	2990	200 ± 21	262 ± 43	261 ± 42	—
Bromoviridae	Ilarvirus	CiVV_RNA1	—	3433	245 ± 17	291 ± 48	287 ± 45	290 ± 42
Bromoviridae	Ilarvirus	CiVV_RNA2	—	2914	227 ± 29	257 ± 41	252 ± 40	—
Bromoviridae	Ilarvirus	ELMV_RNA1	—	3431	195 ± 11	285 ± 46	276 ± 41	279 ± 44
Bromoviridae	Ilarvirus	ELMV_RNA2	—	2874	190 ± 25	254 ± 41	246 ± 43	—
Bromoviridae	Ilarvirus	ParMV_RNA1	—	3518	249 ± 20	292 ± 48	282 ± 44	301 ± 50
Bromoviridae	Ilarvirus	ParMV_RNA2	—	2922	194 ± 21	247 ± 40	248 ± 39	—
Bromoviridae	Ilarvirus	PrDV_RNA1	—	3374	176 ± 24	285 ± 46	273 ± 42	275 ± 47
Bromoviridae	Ilarvirus	PrDV_RNA2	—	2593	149 ± 17	239 ± 39	229 ± 37	232 ± 39
Bromoviridae	Ilarvirus	SpLV_RNA1	—	3439	199 ± 19	291 ± 48	275 ± 44	291 ± 46
Bromoviridae	Ilarvirus	SpLV_RNA2	—	2939	201 ± 22	253 ± 40	237 ± 37	—
Bromoviridae	Ilarvirus	ToSV_RNA1	—	3491	232 ± 28	286 ± 46	286 ± 48	283 ± 46
Bromoviridae	Ilarvirus	ToSV_RNA2	—	2926	202 ± 15	253 ± 42	243 ± 40	—
Bromoviridae	Ilarvirus	TuAMV_RNA1	—	3459	226 ± 17	301 ± 48	292 ± 47	300 ± 48
Bromoviridae	Ilarvirus	TuAMV_RNA2	—	2944	191 ± 9	258 ± 41	246 ± 40	—
Caliciviridae	Nebovirus	caliciNB	—	7453	473 ± 48	502 ± 79	501 ± 77	498 ± 83
Caliciviridae	Nebovirus	newbury	—	7454	372 ± 18	495 ± 78	496 ± 82	498 ± 83
Caliciviridae	Norovirus	murineNorov1	—	7382	380 ± 36	517 ± 81	521 ± 82	491 ± 83
Caliciviridae	Norovirus	norwalk	1ihm	7654	430 ± 35	552 ± 84	551 ± 77	—
Caliciviridae	Sapovirus	porcineSapo	—	7320	361 ± 36	480 ± 77	486 ± 73	—
Caliciviridae	Sapovirus	sapoMc10	—	7458	544 ± 39	486 ± 78	491 ± 73	—
Caliciviridae	Sapovirus	saporo	—	7429	510 ± 33	508 ± 79	509 ± 79	—
Caliciviridae	Vesivirus	rabbitVV	—	8380	401 ± 20	524 ± 81	523 ± 82	—
Caliciviridae	Vesivirus	stellerVV	—	8305	415 ± 16	508 ± 79	521 ± 77	—
Caliciviridae	Vesivirus	VESV	—	8284	374 ± 41	516 ± 76	516 ± 78	—
Comovirinae	Comovirus	BPMV_RNA1	1bmrv	5995	433 ± 40	430 ± 67	434 ± 72	443 ± 66
Comovirinae	Comovirus	BPMV_RNA2	1bmrv	3662	288 ± 26	302 ± 49	298 ± 48	302 ± 50
Comovirinae	Comovirus	CowSMV_RNA1	—	5957	339 ± 28	427 ± 69	425 ± 63	430 ± 62
Comovirinae	Comovirus	CowSMV_RNA2	—	3732	255 ± 30	315 ± 51	302 ± 49	309 ± 54
Comovirinae	Comovirus	CPMV_RNA1	1ny7	5889	360 ± 24	423 ± 66	415 ± 66	439 ± 72
Comovirinae	Comovirus	RadMV_RNA1	—	6064	357 ± 21	427 ± 67	422 ± 63	431 ± 73
Comovirinae	Comovirus	RadMV_RNA2	—	4020	274 ± 20	329 ± 53	315 ± 52	323 ± 50
Comovirinae	Comovirus	RCMV_RNA1	rcmv	6033	396 ± 28	420 ± 65	410 ± 59	417 ± 61
Comovirinae	Comovirus	SquashMV_RNA1	—	5865	474 ± 27	419 ± 69	419 ± 70	436 ± 71
Comovirinae	Comovirus	SquashMV_RNA2	—	3354	226 ± 17	285 ± 48	291 ± 49	288 ± 45
Comovirinae	Comovirus	TurRV_RNA1	—	6082	403 ± 32	440 ± 70	434 ± 70	439 ± 63
Comovirinae	Comovirus	TurRV_RNA2	—	3985	256 ± 18	325 ± 52	304 ± 49	298 ± 46
Comovirinae	Fabavirus	BBWV_RNA1	—	5817	542 ± 37	422 ± 68	428 ± 64	444 ± 74
Comovirinae	Fabavirus	BBWV_RNA2	—	3446	260 ± 24	305 ± 49	303 ± 46	307 ± 49
Comovirinae	Fabavirus	mikaniaMMV_RNA1	—	5862	505 ± 46	433 ± 69	450 ± 67	443 ± 73

Comovirinae	Fabavirus	mikaniaMMV_RNA2	—	3418	259 ± 30	303 ± 49	289 ± 47	285 ± 51
Comovirinae	Fabavirus	patchMMV_RNA1	—	5956	539 ± 24	440 ± 70	428 ± 68	438 ± 62
Comovirinae	Fabavirus	patchMMV_RNA2	—	3591	262 ± 32	320 ± 51	313 ± 51	316 ± 55
Comovirinae	Nepovirus	arabisMV_RNA1	—	7334	318 ± 30	485 ± 74	475 ± 78	468 ± 73
Comovirinae	Nepovirus	arabisMV_RNA2	—	3820	207 ± 20	323 ± 53	307 ± 47	319 ± 45
Comovirinae	Nepovirus	blackCRV_RNA1	—	7711	306 ± 31	502 ± 75	486 ± 76	488 ± 78
Comovirinae	Nepovirus	blackCRV_RNA2	—	6405	315 ± 22	445 ± 67	418 ± 65	417 ± 71
Comovirinae	Nepovirus	raspRV_RNA1	—	7935	404 ± 39	528 ± 83	520 ± 81	539 ± 81
Comovirinae	Nepovirus	raspRV_RNA2	—	3914	203 ± 13	318 ± 50	317 ± 53	319 ± 48
Comovirinae	Nepovirus	TRSV_RNA2	1a6c	7271	257 ± 20	500 ± 80	484 ± 75	502 ± 84
Flaviviridae	Flavivirus	alkhurma	—	10685	714 ± 36	659 ± 10	651 ± 10	646 ± 96
Flaviviridae	Flavivirus	apoi	—	10116	484 ± 38	612 ± 96	626 ± 97	615 ± 88
Flaviviridae	Flavivirus	dengue	—	10735	589 ± 45	654 ± 99	634 ± 98	586 ± 96
Flaviviridae	Flavivirus	montana	—	10690	588 ± 34	649 ± 99	652 ± 10	628 ± 92
Flaviviridae	Flavivirus	powassan	—	10839	674 ± 52	668 ± 10	663 ± 97	656 ± 95
Flaviviridae	Flavivirus	rioBravo	—	10140	520 ± 79	631 ± 10	635 ± 98	618 ± 91
Flaviviridae	Hepacivirus	HepC2	—	9711	443 ± 36	617 ± 10	595 ± 86	604 ± 96
Flaviviridae	Hepacivirus	HepC5	—	9343	538 ± 88	585 ± 97	586 ± 91	580 ± 86
Flaviviridae	Hepacivirus	HepC6	—	9628	440 ± 25	601 ± 92	607 ± 93	599 ± 95
Flaviviridae	Pestivirus	border	—	12333	560 ± 38	681 ± 10	688 ± 10	—
Flaviviridae	Pestivirus	BVDV1	—	12573	547 ± 30	692 ± 10	684 ± 10	—
Flaviviridae	Pestivirus	classicalSFV	—	12301	617 ± 61	675 ± 10	667 ± 98	—
Flaviviridae	Pestivirus	pestiGiraffe	—	12602	598 ± 70	693 ± 11	684 ± 10	—
Picornaviridae	Aphthovirus	BovRBV	—	7556	375 ± 36	486 ± 76	474 ± 77	444 ± 68
Picornaviridae	Aphthovirus	ERAV	2wff	7734	430 ± 33	508 ± 81	483 ± 75	518 ± 82
Picornaviridae	Aphthovirus	FMDV_type0	1zba	8134	406 ± 31	521 ± 81	517 ± 79	504 ± 86
Picornaviridae	Cardiovirus	saffold	—	8115	487 ± 36	523 ± 82	504 ± 78	485 ± 73
Picornaviridae	Cardiovirus	TMEVlike	—	7961	539 ± 37	513 ± 82	512 ± 84	494 ± 76
Picornaviridae	Enterovirus	BEV	1bev	7414	462 ± 47	497 ± 79	484 ± 76	495 ± 88
Picornaviridae	Enterovirus	Hentero107	—	7423	539 ± 31	487 ± 77	480 ± 77	474 ± 71
Picornaviridae	Enterovirus	Hrhino14	1d3i	7212	419 ± 17	458 ± 73	449 ± 71	437 ± 63
Picornaviridae	Erbovirus	ERBV1	—	8828	393 ± 27	548 ± 90	549 ± 86	538 ± 80
Picornaviridae	Kobuvirus	aichi	—	8251	235 ± 20	508 ± 79	491 ± 78	421 ± 63
Picornaviridae	Kobuvirus	bovineKV	—	8374	405 ± 28	533 ± 82	539 ± 79	470 ± 66
Picornaviridae	Kobuvirus	porcineKV	—	8210	266 ± 25	516 ± 79	499 ± 80	445 ± 75
Picornaviridae	Parechovirus	ljungan	—	7590	425 ± 36	490 ± 77	473 ± 75	478 ± 73
Picornaviridae	Sapelovirus	asapelo	—	8289	433 ± 38	520 ± 82	506 ± 77	506 ± 75
Picornaviridae	Senecavirus	SVV	3cji	7310	364 ± 24	480 ± 76	475 ± 76	454 ± 72
Picornaviridae	Teschovirus	ptescho1	—	7117	297 ± 23	482 ± 78	480 ± 73	498 ± 84
Picornaviridae	Tremovirus	AEV	—	7055	425 ± 43	487 ± 77	469 ± 74	488 ± 74
Tombusviridae	Aureusvirus	MaWLMV	—	4293	350 ± 15	357 ± 56	356 ± 55	—
Tombusviridae	Aureusvirus	pothos	—	4354	358 ± 18	361 ± 57	358 ± 56	—
Tombusviridae	Avenavirus	OCSV	—	4114	327 ± 18	331 ± 50	324 ± 51	—
Tombusviridae	Carmovirus	angelonia	—	3964	338 ± 16	322 ± 52	319 ± 49	—
Tombusviridae	Carmovirus	JapINRV	—	4014	326 ± 45	331 ± 52	327 ± 55	—
Tombusviridae	Carmovirus	PelFBV	—	3923	266 ± 13	336 ± 53	327 ± 52	—
Tombusviridae	Carmovirus	TuCrV	3zx8	4050	332 ± 25	333 ± 54	326 ± 53	—
Tombusviridae	Necrovirus	TNV_A	1tnv	3684	269 ± 6	298 ± 47	296 ± 49	—
Tombusviridae	Tombusvirus	GrALV	—	4731	291 ± 22	375 ± 59	384 ± 60	—
Tombusviridae	Tombusvirus	pearLV	—	4766	261 ± 12	367 ± 59	369 ± 57	—
Tymoviridae	Maculavirus	GFkV	—	7564	250 ± 20	440 ± 69	457 ± 69	—
Tymoviridae	Marafivirus	GVSV1	—	6506	392 ± 36	446 ± 68	437 ± 67	—
Tymoviridae	Marafivirus	MRFV	—	6305	451 ± 23	443 ± 67	450 ± 71	—
Tymoviridae	Marafivirus	OBdv	—	6509	328 ± 35	432 ± 72	424 ± 61	—
Tymoviridae	Marafivirus	OLV3	—	7148	312 ± 27	429 ± 67	426 ± 68	—
Tymoviridae	Tymovirus	AnVYV	—	6151	250 ± 17	356 ± 56	357 ± 55	—
Tymoviridae	Tymovirus	ChYMV	—	6517	217 ± 16	357 ± 58	339 ± 53	—
Tymoviridae	Tymovirus	DiYMV	—	6290	223 ± 26	361 ± 56	353 ± 55	—
Tymoviridae	Tymovirus	DuMV	—	6181	336 ± 44	384 ± 60	384 ± 59	—
Tymoviridae	Tymovirus	EgMV	—	6331	186 ± 18	352 ± 57	346 ± 53	—
Tymoviridae	Tymovirus	ErLV	—	6035	248 ± 24	373 ± 60	368 ± 59	—
Tymoviridae	Tymovirus	NeRNV	—	6285	302 ± 23	361 ± 56	351 ± 53	—
Tymoviridae	Tymovirus	OkMV	—	6223	188 ± 29	333 ± 52	333 ± 50	—
Tymoviridae	Tymovirus	OnYMV	—	6211	247 ± 31	384 ± 62	373 ± 57	—

Tymoviridae	Tymovirus	P1MV	–	6154	$369 \pm 26$	$393 \pm 61$	$389 \pm 64$	–
Tymoviridae	Tymovirus	ScMV	–	6206	$217 \pm 18$	$348 \pm 54$	$343 \pm 53$	–

SI Table I: **Set Viral genomes used in this study, including genome length and average MLD values.**  $\langle MLD_{WT} \rangle$  refers to thermal average of the MLD obtained on WT sequences.  $\overline{\langle MLD_{mut} \rangle}$ ,  $\overline{\langle MLD_{UTRs} \rangle}$ ,  $\overline{\langle MLD_{CB} \rangle}$ , refer to average MLD values obtained on synonymously mutated sequences, synonymously mutated sequences with preserved UTRs, and synonymously mutated sequences with preserved UTRs and codon bias, respectively (see Material and Methods in the main text); in these cases an additional averaging over a wide set of possible mutations is performed. Errors are reported as standard deviations.